

**BGGN 213**  
**Hands-on Lab Session**  
 Class 03  
 Barry Grant  
 UC San Diego  
<http://thegrantlab.org/bgg213>

**Class 3: Hands-on section**  
<http://thegrantlab.org/bgg213/>  
 Week 2

Screenshot of the course schedule page:

Week	Date	Topic
2	Fri 10/01/21	Homology, Sequence similarity, Local and global alignment, classic Needleman-Wunsch, Smith-Waterman and BLAST heuristic approaches, Hands on with dot plots, Needleman-Wunsch and BLAST algorithms highlighting their utility and limitations.
3	Wed 10/06/21	<b>Project: Find a gene project assignment</b> (Part 1) Principles of database searching, due in 2 weeks. (Part 2) Sequence analysis, structure analysis and general data analysis with R due at the end of the quarter.
*	Wed 10/06/21	<b>Optional: Advanced sequence alignment and database searching</b> Detecting remote sequence similarity, Database searching beyond BLAST, Substitution matrices, Using PSI-BLAST, Profiles and HMMs, Protein structure comparisons as a gold standard.
4	Fri 10/08/21	Bioinformatics data analysis with R Why do we use R for bioinformatics? R language basics and the RStudio IDE, Major R data structures and functions, Using R interactively from the RStudio console. Introducing Rmarkdown documents.

## Find-a-Gene Project Assignment

- A total of 35% of the course grade will be assigned based on the "find-a-gene project assignment"

## Find-a-Gene Project Assignment

- A total of 35% of the course grade will be assigned based on the "find-a-gene project assignment"
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

## Find-a-Gene Project Assignment

- A total of 35% of the course grade will be assigned based on the “[find-a-gene project assignment](#)”
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the [example report](#) for format and content guidance.

## Find-a-Gene Project Assignment

- A total of 35% of the course grade will be assigned based on the “[find-a-gene project assignment](#)”
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the [example report](#) for format and content guidance.
  - Your responses to questions Q1-Q4 are due 12pm San Diego time on Tuesday **Oct 19th** (10/19/21).
  - The complete assignment, including responses to **all questions**, is due 12pm San Diego time on **Dec 2nd** (12/02/21).

## Find-a-Gene Project Assignment

- A total of 35% of the course grade will be assigned based on the “[find-a-gene project assignment](#)”
- The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.
- You may wish to consult the scoring rubric (in the linked project description) and the [example report](#) for format and content guidance.
  - Your responses to questions Q1-Q4 are due 12pm San Diego time on Tuesday **Oct 19th** (10/19/21).
  - The complete assignment, including responses to **all questions**, is due 12pm San Diego time on **Dec 2nd** (12/02/21).

### Questions:

**[Q1]** Tell me the name of a protein you are interested in. Include the species and the accession number. This can be a human protein or a protein from any other species as long as its function is known.

If you do not have a favorite protein, select human RBP4 or KIF11. Do not use beta globin as this is in the worked example report that I provide you with online.

**[Q2]** Perform a BLAST search against a DNA database, such as a database consisting of genomic DNA or ESTs. The BLAST server can be at NCBI or elsewhere. Include details of the BLAST method used, database searched and any limits applied (e.g. Organism).

Also include the output of that BLAST search in your document. If appropriate, change the font to `courier size 10` so that the results are displayed neatly. You can also screen capture a BLAST output (e.g. all print screen on a PC or on a MAC press `⌘-shift-4`. The pointer becomes a bulls-eye. Select the area you wish to capture and release. The image is saved as a file called `screen_shot_1.png` in your Desktop directory. It is **not** necessary to print out all of the blast results if there are many pages.

On the BLAST results, clearly indicate a match that represents a protein sequence, encoded from some DNA sequence, that is homologous to your query protein. I need to be able to inspect the pairwise alignment you have selected, including the E value and score. It should be labeled a “genomic clone” or “mRNA sequence”, etc. – but include no functional annotation.

In general, [Q2] is the most difficult for students because it requires you to have a “feel” for how to interpret BLAST results. You need to distinguish between a perfect match to your query (i.e. a sequence that is not “novel”), a near match (something that might be “novel”, depending on the results of [Q4]), and a non-homologous result.

If you are having trouble finding a novel gene try restricting your search to an organism that is poorly annotated.

**[Q3]** Gather information about this “novel” protein. At a minimum, show me the protein sequence of the “novel” protein as displayed in your BLAST results from [Q2] as FASTA format (you can copy and paste the aligned sequence subject lines from your BLAST result page if necessary) or translate your novel DNA sequence using a tool called EMBOSS: Transeq at the EBI. Don’t forget to translate all six reading frames: the ORF (open reading frame) is likely to be the longest sequence without a stop codon. It may not start with a methionine if you don’t have the complete coding region. Make sure the sequence you provide includes a header/subject line and is in traditional FASTA format.

Here, tell me the name of the novel protein, and the species from which it derives. It is very unlikely (but still definitely possible) that you will find a novel gene from an organism such as *S. cerevisiae*, human or mouse, because those genomes have already been thoroughly annotated. It is more likely that you will discover a new gene in a genome that is currently being sequenced, such as bacteria or plants or protozoa.

**[Q4]** Prove that this gene, and its corresponding protein, are novel. For the purposes of this project, “novel” is defined as follows. Take the protein sequence (your answer to [Q3]), and use it as a query in a blastp search of the nr database at NCBI.

- If there is a match with 100% amino acid identity to a protein in the database, from the same species, then your protein is NOT novel (even if the match is to a protein with a name such as “unknown”). Someone has already found and annotated this sequence, and assigned it an accession number.
- If the top match reported has less than 100% identity, then it is likely that your protein is novel, and you have succeeded.
- If there is a match with 100% identity, but to a different species than the one you started with, then you have likely succeeded in finding a novel gene.
- If there are no database matches to the original query from [Q1], this indicates that you have partially succeeded: yes, you may have found a new gene, but no, it is not actually homologous to the original query. You should probably start over.

**[Q5]** Generate a multiple sequence alignment with your novel protein, your original query protein, and a group of other members of this family from different species. A typical number of proteins to use in a multiple sequence alignment for this assignment purpose is a minimum of 5 and a maximum of 20 – although the exact number is up to you. Include the multiple sequence alignment in your report. Use Courier font with a size appropriate to fit page width.

Side-note: Indicate your sequence in the alignment by choosing an appropriate name for each sequence in the input unaligned sequence file (i.e. edit the sequence file so that the species, or short common, names (rather than accession numbers) display in the output alignment and in the subsequent answers below). The goal in this step is to create an interesting alignment for building a phylogenetic tree that illustrates species divergence.

**UC San Diego**  
**BGGN 213**  
A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD.

**(Project:) Find a Gene Assignment Part 1**

The **find-a-gene project** is a required assignment for BGGN-213. The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

You may wish to consult the **scoring rubric** at the end of the above linked project description and the **example report** for format and content guidance.

- Your responses to questions Q1-Q4 are due **Wednesday Oct 20th** (10/20/21) at 12pm San Diego time.
- The complete assignment, including responses to all questions, is due **Friday Dec 3rd** (12/03/21) at 12pm San Diego time.
- In both instances your PDF format report should be submitted to GradeScope. Late responses will not be accepted under any circumstances.

**Videos:**

- 3.1 - Project introduction Please note: due dates may differ from those in video.

**UC San Diego**  
**BGGN 213**  
A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD.

**(Project:) Find a Gene Assignment Part 1**

The **find-a-gene project** is a required assignment for BGGN-213. The objective with this assignment is for you to demonstrate your grasp of database searching, sequence analysis, structure analysis and the R environment that we have covered to date in class.

- Your responses to questions **Q1-Q4** are due 12pm San Diego time on Tuesday **Oct 19th** (11/19/21).
- The complete assignment, including responses to **all questions**, is due 12pm San Diego time on Friday **Dec 2nd** (12/02/21).

# Class 3: Hands-on section

<http://thegrantlab.org/bggn213/>

**Class 03**

**UC San Diego**  
**BGGN 213**  
A hands-on introduction to the computer-based analysis of genomic and biomolecular data from the Division of Biological Sciences, UCSD.

**Schedule**

2	Fri	10/01/21	Homology, Sequence similarity, Local and global alignment, classic Needleman-Wunsch, Smith-Waterman and BLAST heuristic approaches, Hands on with dot plots, Needleman-Wunsch and BLAST algorithms highlighting their utility and limitations.
3	Wed	10/06/21	<b>Project: Find a gene project assignment</b> (Part 1) Principles of database searching, due in 2 weeks. (Part 2) Sequence analysis, structure analysis and general data analysis with R due at the end of the quarter.
*	Wed	10/06/21	<b>Optional: Advanced sequence alignment and database searching</b> Detecting remote sequence similarity, Database searching beyond BLAST, Substitution matrices, Using PSI-BLAST, Profiles and HMMs, Protein structure comparisons as a gold standard.
			<b>Bioinformatics data analysis with R</b> Why do we use R for bioinformatics? R language basics and the RStudio IDE, Major R data structures and functions, Using R interactively from the RStudio console. Introducing Rmarkdown documents.
4	Fri	10/08/21	

**► Details:**

Sequence 1: GATTAC  
Sequence 2: GTCGAGGC

Match Score: -4  
Mismatch Score: -1  
Gap Score: -2

Sequence alignment: G T C G A C G C / G A T T A C . .

	G	T	C	G	A	C	G	C	
G	0	-2	-4	-6	-8	-10	-12	-14	-16
A	-2	1	-1	-3	-5	-6	-4	-10	-8
T	-4	-1	0	-2	-4	-3	-1	-9	-7
C	-6	-3	0	-1	-3	-2	-4	-8	-6
G	-8	-5	-2	-1	-2	-4	-6	-8	-7
A	-10	-7	-4	-3	-2	-1	-3	-5	-7
C	-12	-9	-6	-3	-4	-3	0	-2	-4

**► Reference:**  
See the lecture and hands-on session for class 2 for a full discussion of Global, Local, and various Heuristic approaches to biomolecular sequence alignment.  
Barry J Grant.

**R Shiny App**

**NW App Link**

## YOUR TURN!

- There are **four required** and **one optional** hands-on sections including:

- Limits of using BLAST [~10 mins]
- Using PSI-BLAST [~30 mins]
- Examining conservation patterns [~20 mins]
- BREAK [15 mins] —
- [Optional] Using HMMER [~10 mins]
- Divergence of protein sequence and structure [~25 mins]

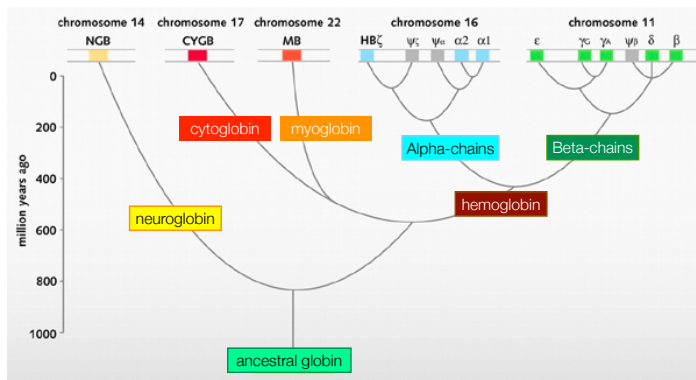
- ▶ Please do answer the last review question (Q20).
- ▶ We encourage discussion at your **Table** and on **Piazza!**

## YOUR TURN!

- There are **four required** and **one optional** hands-on sections including:

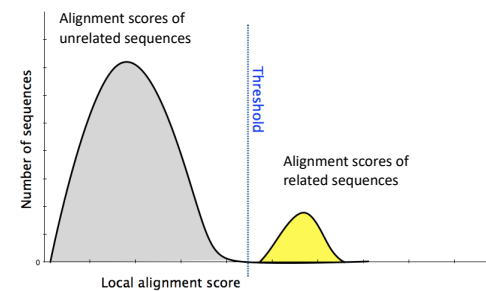
- Limits of using BLAST [~10 mins]
- Using PSI-BLAST [~30 mins]
- Examining conservation patterns [~20 mins]
- BREAK [15 mins] —
- [Optional] Using HMMER [~10 mins]
- Divergence of protein sequence and structure [~25 mins]

- ▶ Please do answer the last review question (Q20).
- ▶ We encourage discussion at your **Table** and on **Piazza!**

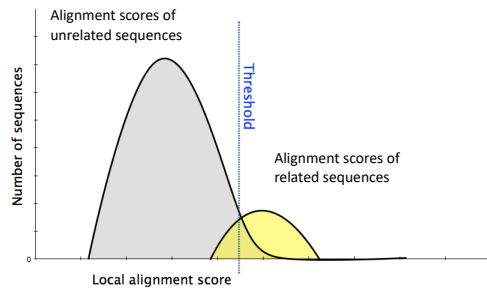


**An evolutionary model of human globins.**  
The different locations of globin genes in human chromosomes are reported at the top of the figure, distinguishing between the functional genes (in color) and the pseudogenes (in grey).

- Ideally, a threshold separates all query related sequences (yellow) from all unrelated sequences (gray)

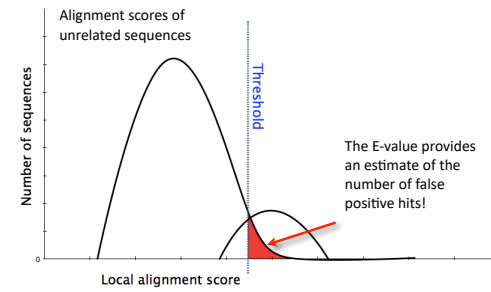


- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



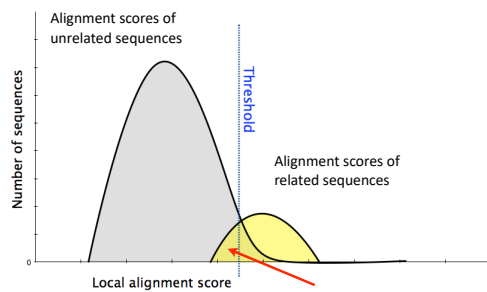
17

- Unfortunately, often both score distributions overlap
  - The E value describes the expected number of hits with a score above the threshold if the query and database are unrelated



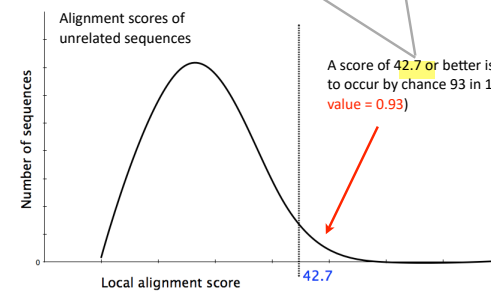
18

- Maybe myoglobin, cytoglobin, neuroglobin etc. are found but not reported because of our E-value cutoff?
  - Lets change the cutoff and see...



19

Description	Max score	Query cover	E value	Max ident	Accession
hemoglobin subunit beta	284	100%	0	100%	NP_000510.1
hemoglobin subunit delta	240	100%	0	75.5%	NP_005321.1
hemoglobin subunit alpha	114	97%	0	43.45%	NP_000508.1
probable ATP-dependent RNA helicase	42.7	10%	0.93	32%	XP_011530405.1



20

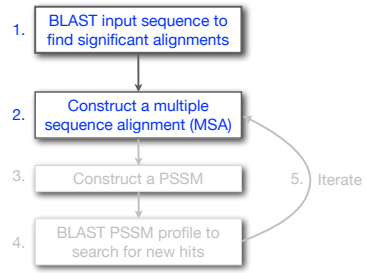
E value: Th alignments with a part





## PSI-BLAST: Position-Specific Iterated BLAST

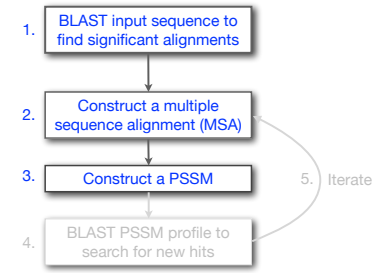
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)

## PSI-BLAST: Position-Specific Iterated BLAST

Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)

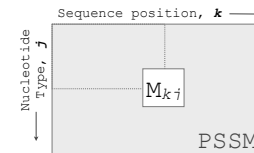
## What is a **PSSM**?

## What are PSSM sequence profiles?

A sequence profile is a **position-specific scoring matrix** (or **PSSM**, often pronounced 'possum') that gives a *quantitative* description of a set of aligned sequences.

PSSMs assign a score to a query sequence and are widely used for database searching.

A simple PSSM has as many columns as there are positions in the alignment, and either 4 rows (one for each DNA nucleotide) or 20 rows (one for each amino acid).



$$M_{kj} = \log \left( \frac{P_{kj}}{P_j} \right)$$

$M_{kj}$  score for the  $j$ th nucleotide at position  $k$   
 $P_{kj}$  probability of nucleotide  $j$  at position  $k$   
 $P_j$  "background" probability of nucleotide  $j$

See Gibskov *et al.* (1987) PNAS 84, 4355



**Example:** Computing a transcription factor bind site PSSM

CCAAATTAGGAAA  
 CCTATTAAGAAAA  
 CCAAATTAGGAAA  
 CCAAATTCGGATA  
 CCCATTTGAAAA  
 CCTATTTAGTATA  
 CCAAATTAGGAAA  
 CCAAATTGGCAAA  
 TCTATTTTGGAAA  
 CCAATTTTCAAAA

Here we have **10 aligned** transcription factor binding site nucleotide sequences

That span **13 positions** (i.e. columns of nucleotides).

We will build a **13 x 4 PSSM** ( $k=13, j=4$ ).

Computing a transcription factor bind site PSSM

CCAAATTAGGAAA  
 CCTATTAAGAAAA  
 CCAAATTAGGAAA  
 CCAAATTCGGATA  
 CCCATTTGAAAA  
 CCTATTTAGTATA  
 CCAAATTAGGAAA  
 CCAAATTGGCAAA  
 TCTATTTTGGAAA  
 CCAATTTTCAAAA

First we will build an alignment **Counts matrix**

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:													
C:													
G:													
T:													

Computing a transcription factor bind site PSSM

CCAAATTAGGAAA  
 CCTATTAAGAAAA  
 CCAAATTAGGAAA  
 CCAAATTCGGATA  
 CCCATTTGAAAA  
 CCTATTTAGTATA  
 CCAAATTAGGAAA  
 CCAAATTGGCAAA  
 TCTATTTTGGAAA  
 CCAATTTTCAAAA

Alignment Counts matrix:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:													
C:													
G:													
T:													

Position k = 1

Computing a transcription factor bind site PSSM

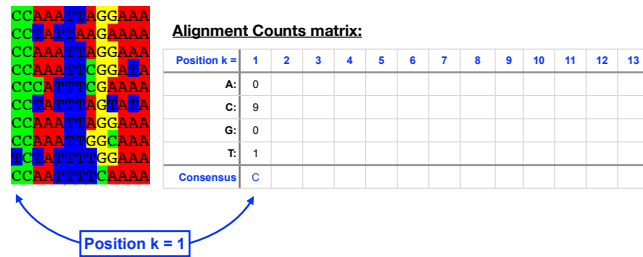
CCAAATTAGGAAA  
 CCTATTAAGAAAA  
 CCAAATTAGGAAA  
 CCAAATTCGGATA  
 CCCATTTGAAAA  
 CCTATTTAGTATA  
 CCAAATTAGGAAA  
 CCAAATTGGCAAA  
 TCTATTTTGGAAA  
 CCAATTTTCAAAA

Alignment Counts matrix:

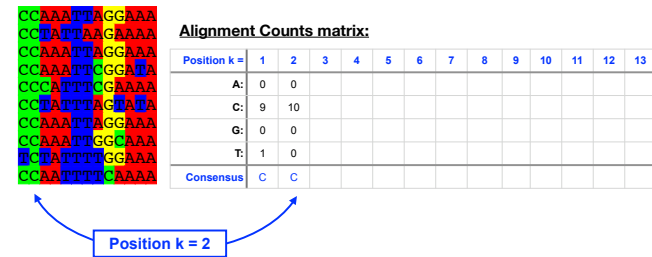
Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0												
C:	9												
G:	0												
T:	1												

Position k = 1

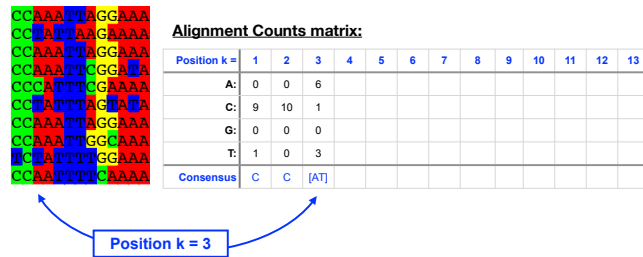
### Computing a transcription factor bind site PSSM



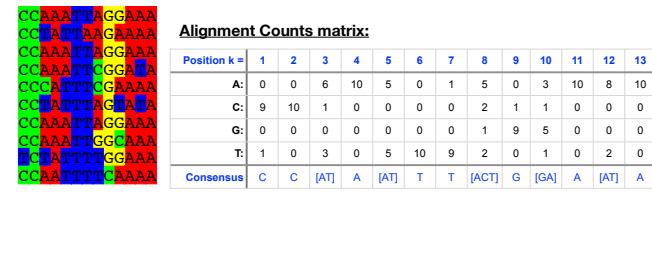
### Computing a transcription factor bind site PSSM



### Computing a transcription factor bind site PSSM



### Computing a transcription factor bind site PSSM



## Computing a transcription factor bind site PSSM

CCAAAATTGGAAA  
 CCTAATTAGGAAAA  
 CCAAAATTAGGAAA  
 CCAAAATCGGATA  
 CCCAATTGGAAAA  
 CCTAATTAGGATA  
 CCAAAATTAGGAAA  
 CCAAAATGGCAAAA  
 CCTAATTGGAAAA  
 CCAAAATTCAAAA

### Alignment Counts matrix:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	C	C	[AT]	A	[AT]	T	T	[ACT]	G	[GA]	A	[AT]	A

Often we will not communicate with the count matrix but rather the derived **average profile** (a.k.a. frequency matrix).

### Average Profile (Frequency) matrix:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	0.6	1	0.5	0	0.1	0.5	0	0.3	1	0.8	1
C:	0.9	1	0.1	0	0	0	0	0.2	0.1	0.1	0	0	0
G:	0	0	0	0	0	0	0	0.1	0.9	0.5	0	0	0
T:	0.1	0	0.3	0	0.5	1	0.9	0.2	0	0.1	0	0.2	0
Consensus	C	C	[AT]	A	[AT]	T	T	[ACT]	G	[GA]	A	[AT]	A

## Computing a transcription factor bind site PSSM

CCAAAATTGGAAA  
 CCTAATTAGGAAAA  
 CCAAAATTAGGAAA  
 CCAAAATCGGATA  
 CCCAATTGGAAAA  
 CCTAATTAGGATA  
 CCAAAATTAGGAAA  
 CCAAAATGGCAAAA  
 CCTAATTGGAAAA  
 CCAAAATTCAAAA

### Alignment Counts matrix:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0
Consensus	C	C	[AT]	A	[AT]	T	T	[ACT]	G	[GA]	A	[AT]	A

Or the "score ( $M_{kj}$ ) matrix" = PSSM

$C_{kj}$  Number of  $j$ th type nucleotide at position  $k$

$Z$  Total number of aligned sequences

$p_j$  "background" probability of nucleotide  $j$

$p_{kj}$  probability of nucleotide  $j$  at position  $k$

$$M_{kj} = \log\left(\frac{p_{kj}}{p_j}\right) \quad p_{kj} = \frac{C_{kj} + p_j}{Z + 1}$$

$$M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right)$$

Adapted from Hertz and Stormo, Bioinformatics 15:563-577

## Computing a transcription factor bind site PSSM...

Alignment Matrix:  $C_{kj}$

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	0	0	6	10	5	0	1	5	0	3	10	8	10
C:	9	10	1	0	0	0	0	2	1	1	0	0	0
G:	0	0	0	0	0	0	0	1	9	5	0	0	0
T:	1	0	3	0	5	10	9	2	0	1	0	2	0

$$k=1, j=A: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{0 + 0.25 / 10 + 1}{0.25}\right) = -2.4$$

$$k=1, j=C: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{9 + 0.25 / 10 + 1}{0.25}\right) = 1.2$$

$$k=1, j=T: M_{kj} = \log\left(\frac{C_{kj} + p_j / Z + 1}{p_j}\right) = \log\left(\frac{1 + 0.25 / 10 + 1}{0.25}\right) = -0.8$$

PSSM:  $M_{kj}$

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4	-2.4
T:	-0.8	-2.4	0.2	-2.4	0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4

## Scoring a test sequence

Query Sequence

CCTAATTAGGATA

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4
T:	-0.8	-2.4	0.2	-2.4	0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4

Test seq: C C T A T T T A G G A T A

$$\text{Query Score} = 1.2 + 1.3 + 0.2 + 1.3 + 0.6 + 1.3 + 1.2 + 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3 = 11.9$$

## Scoring a test sequence

Query Sequence  
**CCTAATTTAGGATA**

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4	-2.4
G:	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4	-2.4
T:	-0.8	-2.4	0.2	-2.4	0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4

Test seq: C C T A T T T A G G A T A

$$\begin{aligned} \text{Query Score} &= 1.2 + 1.3 + 0.2 + 1.3 + 0.6 + 1.3 + 1.2 \\ &\quad + 0.6 + 1.2 + 0.6 + 1.3 + -0.2 + 1.3 \\ &= 11.9 \end{aligned}$$

Q. Does the query sequence match the DNA sequence profile?

## Scoring a test sequence...

Query Sequence      Best Possible Sequence  
**CCTAATTTAGGATA**      **CCAATTTAGGAAA**

PSSM:

Position k =	1	2	3	4	5	6	7	8	9	10	11	12	13
A:	-2.4	-2.4	0.8	1.3	0.6	-2.4	-0.8	0.6	-2.4	0.2	1.3	1.1	1.3
C:	1.2	1.3	-0.8	-2.4	-2.4	-2.4	-2.4	-0.2	-0.8	-0.8	-2.4	-2.4	-2.4
G:	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-2.4	-0.8	1.2	0.6	-2.4	-2.4	-2.4
T:	-0.8	-2.4	0.2	-2.4	0.6	1.3	1.2	-0.2	-2.4	-0.8	-2.4	-0.2	-2.4

Max Score: C C A A T T T A G G A A A

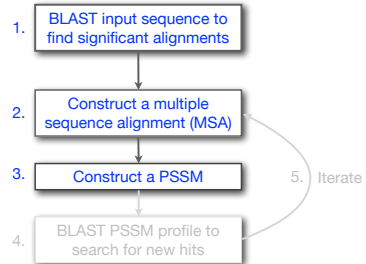
$$\begin{aligned} \text{Max Score} &= 1.2 + 1.3 + 0.8 + 1.3 + 0.6 + 1.3 + 1.2 \\ &\quad + 0.6 + 1.2 + 0.6 + 1.3 + 1.1 + 1.3 \\ &= 13.8 \end{aligned}$$

A. Following method in Harbison *et al.* (2004) Nature 431:99-104

Heuristic threshold for match = 60% x Max Score = (0.6 x 13.8 = 8.28);  
 11.9 > 8.28; Therefore our query is a potential TFBS!

## PSI-BLAST: Position-Specific Iterated BLAST

Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)

## Inspect the blastp output to identify empirical "rules" regarding amino acids tolerated at each position

```

730496 66 FTVDENGQMSATAGRVRLFNNDVDCADHIGSFDTEDPAKFKMRYUGVASFLQKGNDDH 125
200679 63 FSVDEKGHMSATAGRVRLLSNWEVCADHVGTFDTEPAKFKMRYUGVASFLQKGNDDH 122
206589 34 FSVDEKGHMSATAGRVRLLSNWEVCADHVGTFDTEPAKFKMRYUGVASFLQKGNDDH 93
2136812 2 MSATAGRVRLLSNWEVCADHVGTFDTEPAKFKMRYUGVASFLQKGNDDH 53
132408 65 FKIEDNGKTTATAGRVRLDKLELCANNVGTFIETNDPAKFKMRYUGALAILERGLDDH 124
267584 44 FSVDSGRVTTAAGRVIIILNNWECANNFCTFEDTPDPAKFKMRYUGAAAYLQSGNDDH 103
267585 44 FSVDSGRVTTAAGRVIIILNNWECANNFCTFEDTPDPAKFKMRYUGAAAYLQSGNDDH 103
8777608 63 FTIHEDGAMTATAGRVIIILNNWECADHMAFTETTPDPAKFKMRYUGAAAYLQSGNDDH 122
6687453 60 FKVEEDGTHMTATAGRVIIILNNWECANNFCTFEDTEPAKFKMRYUGAAAYLQSGYDDH 119
10697027 61 FKVQEDGTHMTATAGRVIIILNNWECANNFCTFEDTEPAKFKMRYUGAAAYLQSGYDDH 140
13645517 1 HVGTFDTEPAKFKMRYUGVASFLQKGNDDH 32
13925316 38 FSVDSGRKHTATAAGRVIIILNNWECANNFCTFEDTPDPAKFKMRYUGAAAYLQSGNDDH 97
131649 65 YTVVEEDGTHMTASSKGRVRLFGFVVICADHAAQYDPTTPARMYHTYQGLASYLSSGGDNY 126
  
```

M

N,M,L,Y,G

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V	
1 M	-1	-2	0	-2	-1	-2	-2	-1	2	-2	1	2	-2	6	0	-3	-2	-1	-2	-1	
2 K	-1	1	0	1	-4	2	4	-2	0	-3	-3	3	-2	-4	-1	0	-1	-3	-2	-3	
3 W	-3	-3	-4	-5	-3	-2	-3	-3	-1	3	-3	2	4	-3	-2	1	-4	-3	12	2	-3
4 V	0	-3	-3	-4	-1	-3	-3	-4	-1	3	-3	2	4	-3	-2	0	-3	-1	-2	-1	
5 W	-3	-3	-4	-5	-3	-2	-3	-3	-1	3	-3	2	4	-3	-2	1	-4	-3	12	2	-3
6 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0	
7 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1	
8 L	-1	-3	-3	-4	-1	-3	-3	-4	-3	2	2	-3	1	3	-3	-2	-1	-2	0	3	
9 L	-1	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	2	
10 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1	
11 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0	
12 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0	
13 W	-2	-2	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	
14 A	3	2	-1	-2	-1	-2	-1	1	-1	-3	-1	1	-1	-3	-3	-1	1	-3	-3	-1	
15 A	2	2	-1	0	-2	-3	-1	3	0	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	
16 A	4	4	-2	-1	-1	-3	-1	1	0	-3	-2	-1	-1	-3	-1	1	0	-3	-2	-1	
...																					
37 S	2	2	-1	-2	-3	-1	4	1	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	
38 G	0	-3	-1	-2	-3	-2	6	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4	-4	
39 T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-3	-2	0	0	
40 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-2	-1	-4	-3	-3	9	2	-3	-3	
41 Y	-2	-2	-2	-3	-3	-2	-2	-3	2	-2	-1	-2	-1	3	-3	-2	2	7	-1	-1	
42 A	4	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0	

20 amino acids

All the amino acids from position 1 to N (the end of your query protein)

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
1 M	-1	-2	-2	-3	-2	-1	-2	-3	-2	1	2	-2	6	0	-3	-2	-1	-2	-1	1
2 K	-1	1	0	1	-4	2	4	-2	0	-3	-3	3	-2	-4	-1	0	-1	-3	-2	-3
3 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-3	-2	1	-4	-3	12	2	-3	-3
4 V	0	-3	-3	-4	-1	-3	-3	-4	-4	3	1	-3	1	-1	-3	-2	0	-3	-1	4
5 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-3	-2	1	-4	-3	12	2	-3	-3
6 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
7 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1
8 L	-1	-3	-3	-4	-1	-3	-3	-4	-3	2	2	-3	1	3	-3	-2	-1	-2	0	3
9 L	-1	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	2
10 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1
11 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
12 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
13 W	-2	-2	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1
14 A	3	2	-1	-2	-1	-2	-1	1	-1	-3	-1	1	-1	-3	-3	-1	1	-3	-3	-1
15 A	2	2	-1	0	-2	-3	-1	3	0	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2
16 A	4	4	-2	-1	-1	-3	-1	1	0	-3	-2	-1	-1	-3	-1	1	0	-3	-2	-1
...																				
37 S	2	-1	0	-1	-2	-3	-1	4	1	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2
38 G	0	-3	-1	-2	-3	-2	6	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4	-4
39 T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-3	-2	0	0
40 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-2	-1	-4	-3	-3	9	2	-3	-3
41 Y	-2	-2	-2	-3	-3	-2	-2	-3	2	-2	-1	-2	-1	3	-3	-2	2	7	-1	-1
42 A	4	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0

Note: A given amino acid (such as alanine) in your query protein can receive different scores for matching alanine depending on the position in the protein (BLOSUM S<sub>AA</sub> = +4)

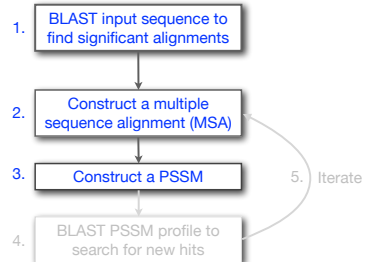
	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
1 M	-1	-2	0	-2	-1	-2	-2	-1	2	-2	1	2	-2	6	0	-3	-2	-1	-2	-1
2 K	-1	1	0	1	-4	2	4	-2	0	-3	-3	3	-2	-4	-1	0	-1	-3	-2	-3
3 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-3	-2	1	-4	-3	12	2	-3	-3
4 V	0	-3	-3	-4	-1	-3	-3	-4	-4	3	1	-3	1	-1	-3	-2	0	-3	-1	4
5 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-3	-2	1	-4	-3	12	2	-3	-3
6 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
7 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1
8 L	-1	-3	-3	-4	-1	-3	-3	-4	-3	2	2	-3	1	3	-3	-2	-1	-2	0	3
9 L	-1	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	2
10 L	-2	-2	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1	1
11 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
12 A	5	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0
13 W	-2	-2	-3	-4	-4	-1	-2	-3	-4	-3	2	4	-3	2	0	-3	-3	-1	-2	-1
14 A	3	2	-1	-2	-1	-2	-1	1	-1	-3	-1	1	-1	-3	-3	-1	1	-3	-3	-1
15 A	2	2	-1	0	-2	-3	-1	3	0	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2
16 A	4	4	-2	-1	-1	-3	-1	1	0	-3	-2	-1	-1	-3	-1	1	0	-3	-2	-1
...																				
37 S	2	-1	0	-1	-2	-3	-1	4	1	-3	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2
38 G	0	-3	-1	-2	-3	-2	6	-2	-4	-4	-2	-3	-4	-2	0	-2	-3	-3	-4	-4
39 T	0	-1	0	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5	-3	-2	0	0
40 W	-3	-3	-4	-5	-3	-2	-3	-3	-3	-3	-2	-2	-1	-4	-3	-3	9	2	-3	-3
41 Y	-2	-2	-2	-3	-3	-2	-2	-3	2	-2	-1	-2	-1	3	-3	-2	2	7	-1	-1
42 A	4	-2	-2	-2	-1	-1	0	-2	-2	-2	-1	-1	-3	-1	1	0	-3	-2	0	0

The PSI-BLAST PSSM is essentially a query customized scoring matrix that is more sensitive than BLOSUM.

Note: A given amino acid (such as alanine) in your query protein can receive different scores for matching alanine depending on the position in the protein (BLOSUM S<sub>AA</sub> = +4)

### PSI-BLAST: Position-Specific Iterated BLAST

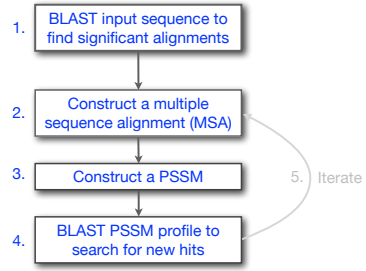
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul et al., Nuc. Acids Res. (1997) 25:3389-3402)

## PSI-BLAST: Position-Specific Iterated BLAST

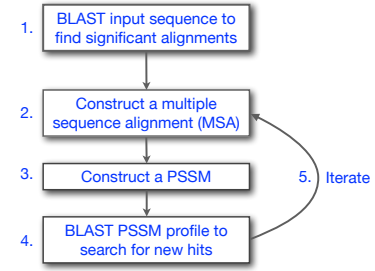
Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



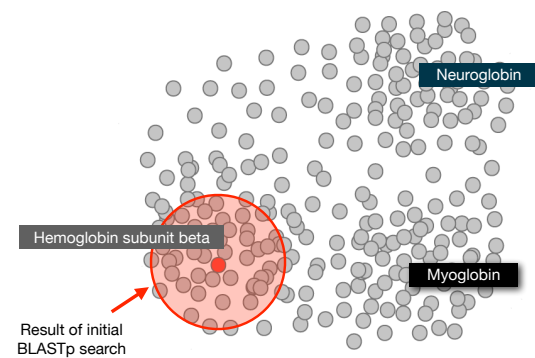
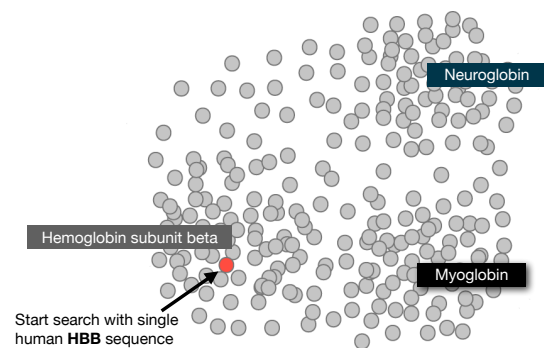
(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)

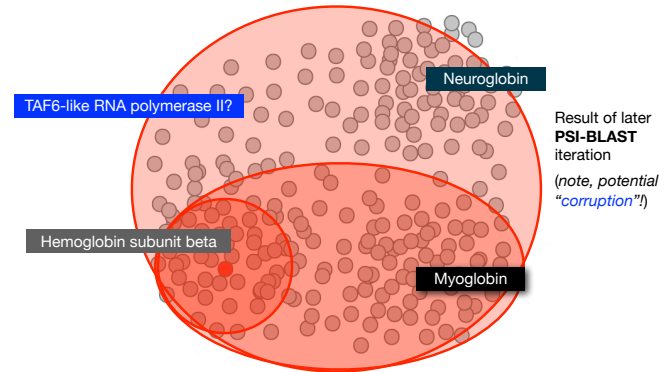
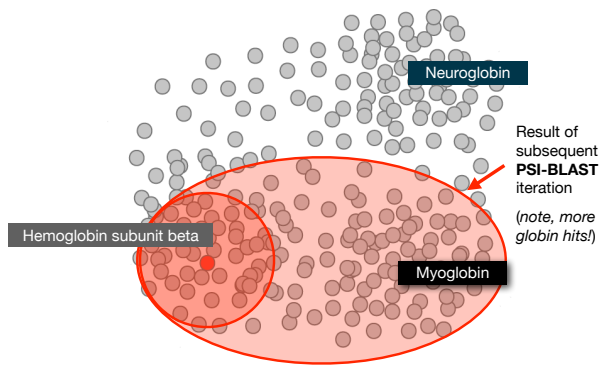
## PSI-BLAST: Position-Specific Iterated BLAST

Many proteins in a database are too distantly related to a query to be detected using standard BLAST. In many other cases matches are detected but are so distant that the inference of homology is unclear. Enter the more sensitive PSI-BLAST



(see Altschul *et al.*, Nuc. Acids Res. (1997) 25:3389-3402)





Description	Max score	Total score	Query cover	E value	Ident	Accession
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1

1

Description	Max score	Total score	Query cover	E value	Ident	Accession
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1
myoglobin [Homo sapiens]	80.5	80.5	97%	2e-19	26%	NP_005359.1
neuroglobin [Homo sapiens]	54.7	54.7	92%	2e-09	23%	NP_067080.1

1

2

New relevant globins found only by PSI-BLAST

Description	Max score	Total score	Query cover	E value	Ident	Accession
hemoglobin subunit beta [Homo sapiens]	301	301	100%	2e-106	100%	NP_000509.1
hemoglobin subunit delta [Homo sapiens]	284	284	100%	7e-100	93%	NP_000510.1
hemoglobin subunit epsilon [Homo sapiens]	240	240	100%	2e-82	76%	NP_005321.1
hemoglobin subunit gamma-2 [Homo sapiens]	235	235	100%	2e-80	73%	NP_000175.1
hemoglobin subunit gamma-1 [Homo sapiens]	232	232	100%	3e-79	73%	NP_000550.2
hemoglobin subunit alpha [Homo sapiens]	114	114	97%	7e-33	43%	NP_000508.1
hemoglobin subunit zeta [Homo sapiens]	100	100	97%	3e-27	36%	NP_005323.1
myoglobin [Homo sapiens]	80.5	80.5	97%	2e-19	26%	NP_005359.1
neuroglobin [Homo sapiens]	54.7	54.7	92%	2e-09	23%	NP_067080.1
myoglobin [Homo sapiens]	159	159	97%	3e-50	26%	NP_005359.1
hemoglobin subunit alpha [Homo sapiens]	151	151	97%	3e-47	42%	NP_000508.1
hemoglobin subunit mu [Homo sapiens]	147	147	97%	6e-46	35%	NP_01003938.1
hemoglobin subunit theta-1 [Homo sapiens]	147	147	97%	2e-45	37%	NP_005322.1
neuroglobin [Homo sapiens]	134	134	92%	3e-40	23%	NP_067080.1
PREDICTED: cytoglobin isoform X2 [Homo sapiens]	115	115	66%	3e-33	25%	XP_016879605.1
PREDICTED: microtubule cross-linking factor 1 isoform X1 [Homo sapiens]	46.3	46.3	27%	7e-06	39%	XP_011523942.1
PREDICTED: microtubule cross-linking factor 1 isoform X4 [Homo sapiens]	46.3	46.3	27%	7e-06	39%	XP_00528156.1

Inclusion of irrelevant hits can lead to PSSM corruption

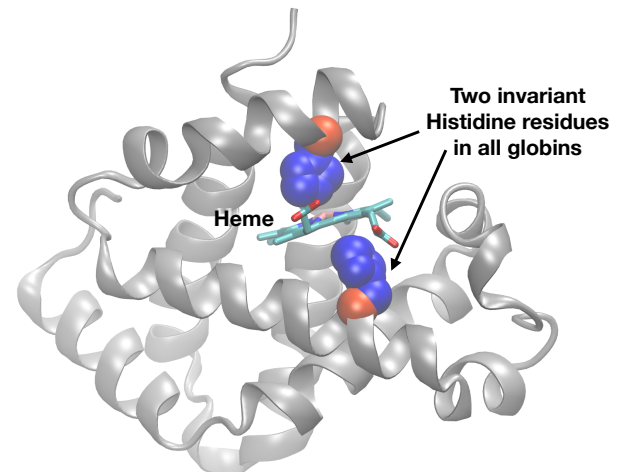
## YOUR TURN!

- There are **four required** and **one optional** hands-on sections including:

- Limits of using BLAST [~10 mins]
- Using PSI-BLAST [~30 mins]
- Examining conservation patterns [~20 mins]  
— BREAK [15 mins] —
- [Optional] Using HMMER [~10 mins]
- Divergence of protein sequence and structure [~25 mins]

- Please do answer the last review question (Q20).
- We encourage discussion at your **Table** and on **Piazza**!

Query_73613	1	MVHLTPEEKSAVTALMGKRV--NVEVGGEGALGRLLVVPWQRFEE--SFGDLSTPDAMV--GNPKVKAHGKVLGAF	72
NP_000510.1	1	MVHLTPEEKSAVTALMGKRV--NVEVGGEGALGRLLVVPWQRFEE--SFGDLSTPDAMV--GNPKVKAHGKVLGAF	72
NP_000175.1	1	MGHFTEDDKATITSLMGKV--NVEDAGSETLGRLLVVPWQRFEE--SFGDLSSASAIM--GNPKVKAHGKVLGAF	72
NP_000509.1	1	MVHLTPEEKSAVTALMGKRV--NVEVGGEGALGRLLVVPWQRFEE--SFGDLSTPDAMV--GNPKVKAHGKVLGAF	72
NP_005321.1	1	MVHFTAEKKAAVTSLMSKM--NVEEAGGEGALGRLLVVPWQRFEE--SFGDLSSASAIM--GNPKVKAHGKVLGAF	72
NP_000550.2	1	MGHFTEDDKATITSLMGKV--NVEDAGSETLGRLLVVPWQRFEE--SFGDLSSASAIM--GNPKVKAHGKVLGAF	72
NP_005322.1	1	--MLTKTERTIIVSMKAKISTQADTIGTETLEELFLSPQTKYTF--HF-----DLHGSAQLRAHGKVVAAV	67
NP_000508.1	1	--MVLSPADKTNVKAAMKVAHAGEYGAELERMFSPQTKYTF--HF-----DLHGSAQLRAHGKVVAAV	67
XP_005257062.1	1	[15]SEELSEAEKKAQAMMARLYANCEDGVALLVRFVYFNSAKQYFS--QFKHMEDELEME--RSPQLRKHACRVMGAL	89
NP_001003938.1	1	--MLSAEQRAIQAQNDLIGAGHEAQGAELLRFLPSTKQYTF--HL-----SACQ--DATQLLHGQRMLAAV	66
NP_005322.1	1	--MALSAEDRALVRLMKKLGSNVGVYVTEALERTFLAFPATKYTF--H-----LDLGSQSVRAHGKVVAAV	67
NP_593030.1	1	[15]SEELSEAEKKAQAMMARLYANCEDGVALLVRFVYFNSAKQYFS--QFKHMEDELEME--RSPQLRKHACRVMGAL	89
XP_016879605.1	1	-----MEDLEME--RSPQLRKHACRVMGAL	24
NP_001349775.1	1	--MGLSDGEWQLVLANWGRVADIPGRGQVLIILRFGHPETLERFP--KFKLKSDEEMK--ASEDLKKGATVITAL	73
NP_067080.1	1	---MERPEPELIRGSRVRSRSPLEGTVLPAFLALEPDLPLFQYRCQFSSPDECL--SSPEFLIRIHWVLI	72
NP_001369741.1	1	-----MK--ASEDLKKGATVITAL	18
Query_73613	73	SDGLAHLNDLKGK---FATLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPPVQAAYQKVVAGVANALAHKYH	147
NP_000510.1	73	SDGLAHLNDLKGK---FQGLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPPVQAAYQKVVAGVANALAHKYH	147
NP_000175.1	73	GDAIKHLDDLKGT---FAQLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPEVQASQWQKVVAGVANALAHKYH	147
NP_000509.1	73	SDGLAHLNDLKGK---FATLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPPVQAAYQKVVAGVANALAHKYH	147
NP_005321.1	73	GDAIKHNDLKPA---FAKLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPEVQASQWQKVVAGVANALAHKYH	147
NP_000550.2	73	GDAIKHLDDLKGT---FAQLSELHCDKLVDPENFRLLGNLVLCVLAHFGKEFTPEVQASQWQKVVAGVANALAHKYH	147
NP_005322.1	68	GDVAKSIDDIGGA---LAKLSELHAYILRVDVPMFKLSHCLLVTLAARFADTAEAHAANDKFLSVSVLTSKRY	142
NP_000509.1	68	TNAVAVHDDMPNA---LALSGLHAKLIRVDVPMFKLSHCLLVTLAARFADTAEAHAANDKFLSVSVLTSKRY	142
XP_005257062.1	90	NTVVENLHDPDKVevLALVQKALAKHKVEPVYFKLSGVILEVVAEEFASDPPETQRAWKLEGLVSVTAAYK [35]	202
NP_001003938.1	67	GAAYQVYDNRRAA---LSPALGLIAKLVDPANFVFLIQCFVYVLAHSHGDFTVVQAAWREKFLVAVVLTETERY	141
NP_005322.1	68	SLAVYKLDLDMFA---LQALSHLHACGLRVDPAFQGLLQCLLIVLANVPCDFVQAAWREKFLVAVVLTETERY	142
NP_593030.1	90	NTVVENLHDPDKVevLALVQKALAKHKVEPVYFKLSGVILEVVAEEFASDPPETQRAWKLEGLVSVTAAYK [23]	190
XP_016879605.1	25	NTVVENLHDPDKVevLALVQKALAKHKVEPVYFKLSGVILEVVAEEFASDPPETQRAWKLEGLVSVTAAYK [35]	137
NP_001349775.1	74	GGILKKGHEAE---IKPLAQSHATKHKIPVKYLEFISECIQVLSQKHPGDFGADQAGMNAKLEFRKDMASNYK [6]	154
NP_067080.1	73	DAAVTVNDELSSLeYLASLGRKRA--VGKLSFSFTVGSLLYMLKXGLPATPATRAAASQGLYGAVVQMSRNGD [2]	151
NP_001369741.1	19	GGILKKGHEAE---IKPLAQSHATKHKIPVKYLEFISECIQVLSQKHPGDFGADQAGMNAKLEFRKDMASNYK [6]	99





## YOUR TURN!

- There are **four required** and **one optional** hands-on sections including:

- Limits of using BLAST [**~10 mins**]
- Using PSI-BLAST [**~30 mins**]
- Examining conservation patterns [**~20 mins**]
- [Optional] Using HMMER** [**~10 mins**]
- Divergence of protein sequence and structure [**~25 mins**]

— BREAK [15 mins] —

- Please do answer the last review question (**Q20**).
- We encourage discussion at your **Table** and on **Piazza!**

## Problems with PSSMs: Positional dependencies

Do not capture positional dependencies

**WEIRD**  
**WEIRD**  
**WEIQH**  
**WEIRD**  
**WEIQH**

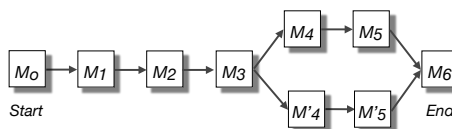
D				0.6
E	I			
H				0.4
I		I		
Q			0.4	
R			0.6	
W	I			

**Note:** We never see **QD** or **RH**, we only see **RD** and **QH**.  
However,  $P(RH)=0.24$ ,  $P(QD)=0.24$ , while  $P(QH)=0.16$

## Markov chains: Positional dependencies ✓

The connectivity or **topology** of a Markov chain can easily be designed to capture dependencies and variable length motifs.

**WEIRD**  
**WEIRD**  
**WEIQH**  
**WEIRD**  
**WEIQH**

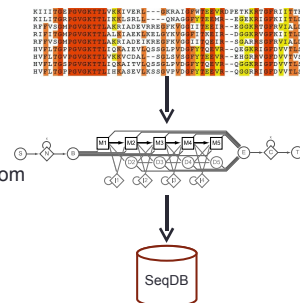


Recall that a PSSM for this motif would give the sequences **WEIRD** and **WEIRH** equally good scores even though the **RH** and **QR** combinations were not observed

## Use of HMMER

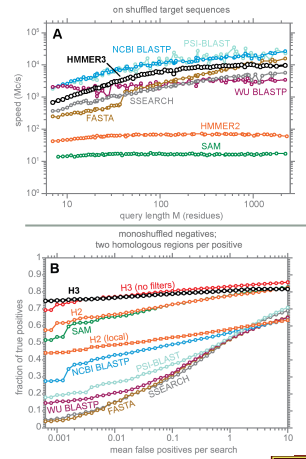
- Widely used by protein family databases

- Use 'seed' alignments
- Until 2010
- Computationally expensive
- Restricted to HMMs constructed from multiple sequence alignments
- Command line application



# HMMER vs BLAST

	HMMER	BLAST
Program	<i>PHMMER</i>	<i>BLASTP</i>
Query	Single sequence	
Target Database	Sequence database	
Program	<i>HMMSCAN</i>	<i>RP-BLAST</i>
Query	Single sequence	
Target Database	Profile HMM database, e.g. Pfam	PSSM database, e.g. CDD
Program	<i>HMMSEARCH</i>	<i>PSS-BLAST</i>
Query	Profile HMM	PSSM
Target Database	Sequence database	
Program	<i>JACKHMMER</i>	<i>PSS-BLAST</i>
Query	Single sequence	
Target Database	Sequence database	

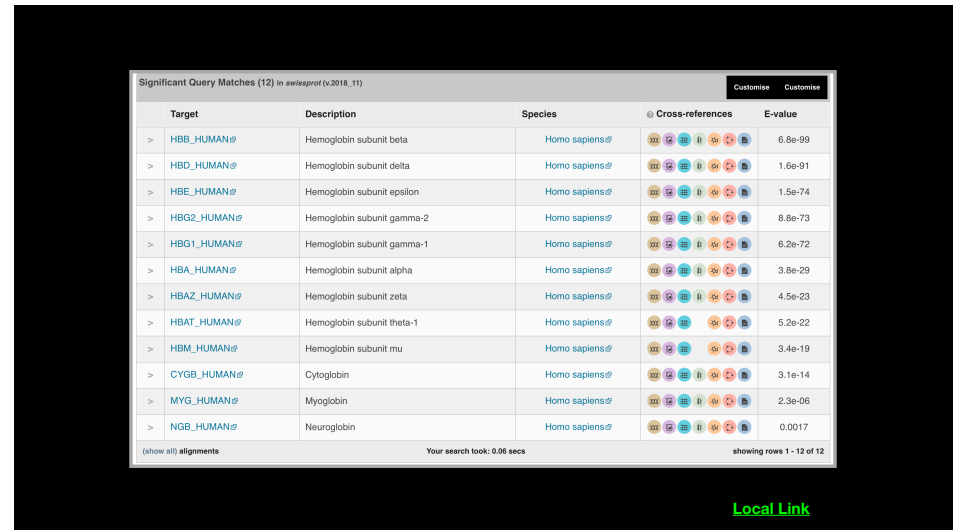
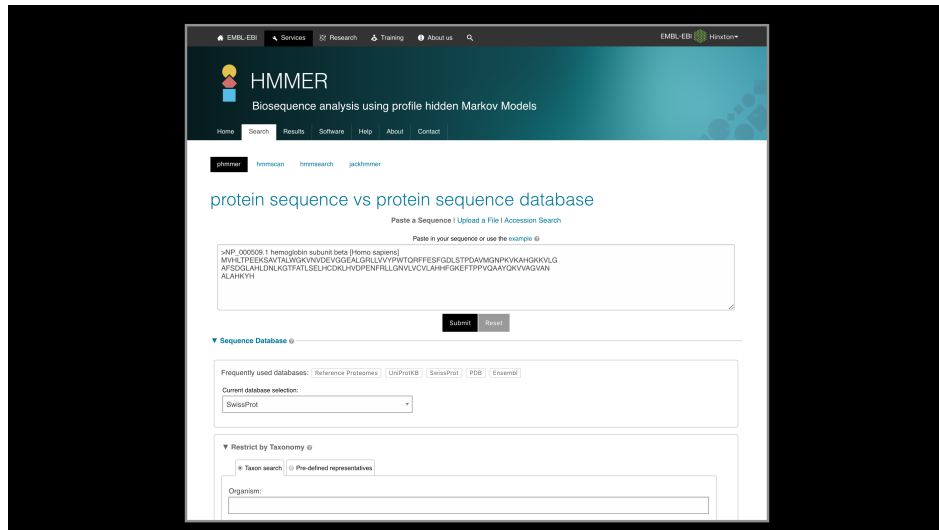


Modified from: S. R. Eddy  
PLoS Comp. Biol., 7:e1002195, 2011.



# Fast Web Searches

- Parallelized searches across compute farm
  - Average query returns ~1 sec
- Range of sequence databases
  - Large Comprehensive
  - Curated / Structure
  - Metagenomics
  - Representative Proteomes
- Family Annotations
  - Pfam
- Batch and RESTful API
  - Automatic and Human interface





## Summary

- **Find a gene project:** You can start working on this now. Submit your responses to Q1-Q4 to get feedback.
- **PSI-BLAST algorithm:** Application of iterative position specific scoring matrices (PSSMs) to improve BLAST sensitivity
- **Hidden Markov models (HMMs):** More versatile probabilistic model for detection of remote similarities
- **Structure comparisons as gold standards:** Structure is more conserved than sequence

## Homework: DataCamp!

Install **R** and **RStudio** (see website)

Complete the **Introduction to R** course on **DataCamp**  
(Check Piazza for your DataCamp invite and sign up with your UCSD email (i.e. first part of your email address) please.)

Let me know **NOW** if you don't have access to DataCamp!